

Efficient Sparseness-Enforcing Projections

Markus Thom¹ and Günther Palm²

Abstract. We propose a linear time and constant space algorithm for computing Euclidean projections onto sets on which a normalized sparseness measure attains a constant value. These non-convex target sets can be characterized as intersections of a simplex and a hypersphere. Some previous methods required the vector to be projected to be sorted, resulting in at least quasilinear time complexity and linear space complexity. We improve on this by adaptation of a linear time algorithm for projecting onto simplexes. In conclusion, we propose an efficient algorithm for computing the product of the gradient of the projection with an arbitrary vector.

1 Introduction

In a great variety of classical machine learning problems, sparse solutions are appealing because they provide more efficient representations compared to non-sparse solutions. Several formal sparseness measures have been proposed in the past and their properties have been thoroughly analyzed [1]. One remarkable sparseness measure is the normalized ratio of the L_1 norm and the L_2 norm of a vector, as originally proposed by [2]:

$$\sigma: \mathbb{R}^n \setminus \{0\} \rightarrow [0, 1], \quad x \mapsto \frac{\sqrt{n} - \frac{\|x\|_1}{\|x\|_2}}{\sqrt{n} - 1}.$$

Here, higher values of σ indicate more sparse vectors. The extreme values of 0 and 1 are achieved for vectors where all entries are equal and vectors where all but one entry vanish, respectively. Further, σ is scale-invariant, that is $\sigma(\alpha x) = \sigma(x)$ for all $\alpha \neq 0$ and all $x \in \mathbb{R}^n \setminus \{0\}$.

The incorporation of explicit sparseness constraints to existing optimization problems while still being able to efficiently compute solutions to them was made possible by [2] through proposition of an operator, which computes the Euclidean projection onto sets on which σ attains a desired value. In other words, given a target degree of sparseness $\sigma^* \in (0, 1)$ with respect to σ , numbers $\lambda_1, \lambda_2 > 0$ can be derived such that $\sigma \equiv \sigma^*$ on the non-convex set

$$D := \{s \in \mathbb{R}_{\geq 0}^n \mid \|s\|_1 = \lambda_1 \text{ and } \|s\|_2 = \lambda_2\}.$$

Clearly, either of λ_1 and λ_2 has to be fixed to a pre-defined value, for example by setting $\lambda_2 := 1$ for achieving normalized vectors, as only their ratio is important in the definition of σ . By restricting possible solutions to certain optimization problems to lie in D , projected gradient descent methods [3] can be used to achieve solutions that fulfill explicit sparseness constraints.

¹driveU / Institute of Measurement, Control and Microtechnology, Ulm University, Ulm, Germany

²Institute of Neural Information Processing, Ulm University, Ulm, Germany

E-mail addresses: markus.thom@uni-ulm.de, guenther.palm@uni-ulm.de

The projection operator of [2] was motivated by geometric ideas, in such that the intersection of hyperplanes, hyperspheres and the non-negative orthant were considered and a procedure of alternating projections was proposed. This procedure is known to produce correct projections when applied to the intersection of convex sets [4]. In the non-convex setup considered here, it is not clear in the first place whether the method also computes correct projections. The results of [5], however, show that alternating projections also work for the sparseness projection, and that the projection onto D is unique almost everywhere.

It was further noted recently that the method of Lagrange multipliers can also be used to derive an implicit, compact representation of the sparseness projection [6]. The algorithm proposed there needs to sort the vector that is to be projected and remember the sorting permutation, resulting in a computational complexity that is quasilinear and a space complexity that is linear in the problem dimensionality n . In this work, we provide a detailed derivation of the results of [6]. Then, by transferring the ideas of [7] to efficiently compute projections onto simplexes we use the implicit representation to propose a linear time and constant space algorithm for computing projections onto D . Ultimately, we propose an algorithm that efficiently computes the product of the gradient of the projection onto D with an arbitrary vector.

2 Notation and Prerequisites

We denote the set of Boolean values with \mathbb{B} , the real numbers with \mathbb{R} and the n -dimensional Euclidean space with \mathbb{R}^n . Subscripts for elements from \mathbb{R}^n denote individual coordinates. All entries of the vector $e \in \mathbb{R}^n$ are unity. $\mathbb{R}^{n \times n}$ is the ring of matrices with n rows and n columns, and $E_n \in \mathbb{R}^{n \times n}$ is the identity matrix. It is well-known that the L_1 norm and the L_2 norm are equivalent in the topological sense [8]:

Remark 1. For all $x \in \mathbb{R}^n$, we have that $\|x\|_2 \leq \|x\|_1 \leq \sqrt{n}\|x\|_2$. If x is sparsely populated, then the latter inequality can be sharpened to $\|x\|_1 \leq \sqrt{d}\|x\|_2$, where $d := \|x\|_0 \leq n$ denotes the number of non-vanishing entries in x .

Therefore $\lambda_2 < \lambda_1 < \sqrt{n}\lambda_2$ must hold for the target norms to achieve a sparseness of $\sigma^* \in (0, 1)$. The projection onto a set contains all points with infimal distance to the projected vector [4]:

Definition 2. Let $x \in \mathbb{R}^n$ and $\emptyset \neq M \subseteq \mathbb{R}^n$. Then every point in

$$\text{proj}_M(x) := \{y \in M \mid \|y - x\|_2 \leq \|z - x\|_2 \text{ for all } z \in M\}$$

is called *Euclidean projection* of x onto M . If there is exactly one point y in $\text{proj}_M(x)$, then $y = \text{proj}_M(x)$ is written for abbreviation.

We further note that projections onto permutation-invariant sets are order-preserving:

Proposition 3. Let $\emptyset \neq M \subseteq \mathbb{R}^n$ such that $P_\tau x \in M$ for all $x \in M$ and all permutation matrices $P_\tau \in \mathbb{R}^{n \times n}$. Let $x \in \mathbb{R}^n$ and $p \in \text{proj}_M(x)$. Then $x_i > x_j$ implies $p_i \geq p_j$ for all $i, j \in \{1, \dots, n\}$.

Throughout the paper, we assume that the input vector to the projection operator is chosen such that the projection onto D is unique. As has been shown by [5], this is fulfilled by almost all $x \in \mathbb{R}^n \setminus \{0\}$ and is thus no restriction in practice.

3 Implicit Representation of the Projection

As noted by [6], the method of Lagrange multipliers can be used to derive an implicit representation of the projection onto D . To make this paper as self-contained as possible, we include an elaborate derivation of their result.

Lemma 4. Let $x \in \mathbb{R}_{\geq 0}^n \setminus D$ such that $\text{proj}_D(x)$ is unique. Then there exist unique numbers $\alpha \in \mathbb{R}$ and $\beta \in \mathbb{R}_{>0}$ such that $\text{proj}_D(x) = \max \left(\frac{1}{\beta} (x - \alpha \cdot e), 0 \right)$.

Proof. We want to find a point $p \in D$ such that the Euclidean distance $\|p - x\|_2$ is minimized. Such a point is guaranteed to exist by the Weierstraß extreme value theorem. The constrained optimization problem leads to the Lagrangian $\mathcal{L}: \mathbb{R}^n \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}_{\geq 0}^n \rightarrow \mathbb{R}$,

$$(p, \alpha, \beta, \gamma) \mapsto \frac{1}{2} \|p - x\|_2^2 + \alpha (\|p\|_1 - \lambda_1) + \frac{\beta-1}{2} (\|p\|_2^2 - \lambda_2^2) - \gamma^T p,$$

where the multiplier β was linearly transformed for notational convenience. By taking the derivative for p and setting it to zero we obtain

$$\frac{\partial \mathcal{L}}{\partial p_i} = \beta p_i - x_i + \alpha - \gamma_i \stackrel{!}{=} 0, \text{ and hence } p_i = \frac{x_i - \alpha + \gamma_i}{\beta} \text{ for all } i \in \{1, \dots, n\}.$$

The complementary slackness condition, $\gamma_i p_i = 0$ for all $i \in \{1, \dots, n\}$, must be satisfied in a local minimum of \mathcal{L} . Hence $p_i > 0$ implies $\gamma_i = 0$ and $p_i = 0$ implies $\gamma_i \geq 0$ for all $i \in \{1, \dots, n\}$. Let $I := \{i \in \{1, \dots, n\} \mid p_i > 0\}$ denote the set of coordinates in which p does not vanish, and let $d := |I|$ denote its cardinality. We have $d \geq 2$, because $d = 0$ is impossible due to $\lambda_1, \lambda_2 > 0$ and $d = 1$ is impossible because $\lambda_1 \neq \lambda_2$. Further, $\gamma_i = 0$ for all $i \in I$ from the complementary slackness condition. Let $\tilde{x} \in \mathbb{R}_{\geq 0}^d$ be the vector with all entries from x with index in I , that is when $I = \{i_1, \dots, i_d\}$ then $\tilde{x}^T = (x_{i_1}, \dots, x_{i_d})$. Note that because all entries of p and x are non-negative, the sum over their entries is identical to their L_1 norm. By taking the derivative of the Lagrangian for α and setting it to zero we have that

$$\lambda_1 = \|p\|_1 = \sum_{i \in I} p_i = \sum_{i \in I} \frac{1}{\beta} (x_i - \alpha) = \frac{1}{\beta} (\|\tilde{x}\|_1 - d\alpha).$$

Analogously, taking the derivative for β and setting it to zero yields

$$\lambda_2^2 = \|p\|_2^2 = \frac{1}{\beta^2} \sum_{i \in I} (x_i - \alpha)^2 = \frac{1}{\beta^2} (\|\tilde{x}\|_2^2 - 2\alpha \|\tilde{x}\|_1 + d\alpha^2).$$

By squaring the expression for λ_1 and dividing by λ_2^2 we get

$$\frac{\lambda_1^2}{\lambda_2^2} = \frac{\|\tilde{x}\|_1^2 - 2d\alpha \|\tilde{x}\|_1 + d^2\alpha^2}{\|\tilde{x}\|_2^2 - 2\alpha \|\tilde{x}\|_1 + d\alpha^2},$$

which leads to the quadratic equation

$$0 = \underbrace{d \left(d - \frac{\lambda_1^2}{\lambda_2^2} \right)}_{=:a} \cdot \alpha^2 + \underbrace{2 \|\tilde{x}\|_1 \left(\frac{\lambda_1^2}{\lambda_2^2} - d \right)}_{=:b} \cdot \alpha + \underbrace{\left(\|\tilde{x}\|_1^2 - \frac{\lambda_1^2}{\lambda_2^2} \|\tilde{x}\|_2^2 \right)}_{=:c}.$$

Before considering the discriminant of this equation, we first note that $d \|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2 \geq 0$ with Remark 1. As p exists by the Weierstraß extreme value theorem and has by definition d non-zero entries, we also have that $d - \frac{\lambda_1^2}{\lambda_2^2} \geq 0$ using Remark 1. Thus we obtain

$$\begin{aligned} D &:= b^2 - 4ac = 4 \|\tilde{x}\|_1^2 \left(d - \frac{\lambda_1^2}{\lambda_2^2}\right)^2 - 4d \left(d - \frac{\lambda_1^2}{\lambda_2^2}\right) \left(\|\tilde{x}\|_1^2 - \frac{\lambda_1^2}{\lambda_2^2} \|\tilde{x}\|_2^2\right) \\ &= 4 \frac{\lambda_1^2}{\lambda_2^2} \left(d - \frac{\lambda_1^2}{\lambda_2^2}\right) \left(d \|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2\right) \geq 0, \end{aligned}$$

so α must be a real number. Solving the equation leads to two possible values for α :

$$\alpha \in \left\{ \frac{-b \pm \sqrt{D}}{2a} \right\} = \left\{ \frac{1}{d} \left(\|\tilde{x}\|_1 \pm \lambda_1 \sqrt{\frac{d \|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2}{d \lambda_2^2 - \lambda_1^2}} \right) \right\}.$$

We first assume that α is the number that arises from the "+" before the square root. From $\lambda_1 = \|p\|_1$ we then obtain

$$\beta = \frac{1}{\lambda_1} (\|\tilde{x}\|_1 - d\alpha) = -\sqrt{\frac{d \|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2}{d \lambda_2^2 - \lambda_1^2}} < 0.$$

With $d \geq 2$ there are two indices $i, j \in I$ with $x_i > x_j$. The derivative of \mathcal{L} for p and the complementary slackness condition then yield $p_i - p_j = \frac{1}{\beta} (x_i - \alpha - x_j + \alpha) = \frac{1}{\beta} (x_i - x_j) < 0$, which contradicts the order-preservation as guaranteed by Proposition 3. Therefore, the choice of α was not correct in the first place, and thus

$$\alpha = \frac{1}{d} \left(\|\tilde{x}\|_1 - \lambda_1 \sqrt{\frac{d \|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2}{d \lambda_2^2 - \lambda_1^2}} \right) \text{ and } \beta = \sqrt{\frac{d \|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2}{d \lambda_2^2 - \lambda_1^2}} > 0$$

must hold. Let $i \in I$, then $0 < p_i = \frac{1}{\beta} (x_i - \alpha)$, and because $\beta > 0$ follows $x_i > \alpha$. For $i \notin I$ it is $0 = p_i = \frac{1}{\beta} (x_i - \alpha + \gamma_i)$ where $\gamma_i \geq 0$, so $0 = x_i - \alpha + \gamma_i \geq x_i - \alpha$, or equivalently $x_i \leq \alpha$. Ultimately, we have that $p_i = \max(\frac{1}{\beta} (x_i - \alpha), 0)$ for all $i \in \{1, \dots, n\}$.

For the claim to hold, it now remains to be shown that α and β are unique. With the uniqueness of the projection p , we thus have to show that from

$$p = \max\left(\frac{1}{\beta_1} (x - \alpha_1 \cdot e), 0\right) = \max\left(\frac{1}{\beta_2} (x - \alpha_2 \cdot e), 0\right)$$

for $\alpha_1, \alpha_2, \beta_1, \beta_2 \in \mathbb{R}$ follows that $\alpha_1 = \alpha_2$ and $\beta_1 = \beta_2$. As shown earlier, there are two distinct indices $i, j \in I$ with $x_i \neq x_j$ and $p_i, p_j > 0$. We hence obtain

$$p_i = \frac{1}{\beta_1} (x_i - \alpha_1) = \frac{1}{\beta_2} (x_i - \alpha_2) \text{ and } p_j = \frac{1}{\beta_1} (x_j - \alpha_1) = \frac{1}{\beta_2} (x_j - \alpha_2),$$

and thus $\frac{p_i}{p_j} = \frac{x_i - \alpha_1}{x_j - \alpha_1} = \frac{x_i - \alpha_2}{x_j - \alpha_2}$. Therefore,

$$0 = (x_i - \alpha_1)(x_j - \alpha_2) - (x_i - \alpha_2)(x_j - \alpha_1) = \alpha_1(x_i - x_j) - \alpha_2(x_i - x_j) = (\alpha_1 - \alpha_2)(x_i - x_j).$$

With $x_i \neq x_j$ we have that $\alpha_1 = \alpha_2$, and substitution in either of p_i or p_j shows that $\beta_1 = \beta_2$. \square

We note that the crucial point in the computation of α is finding the set I where the projection has positive coordinates. With the statement of Proposition 3 the argument of the projection can be sorted before-hand such that $I = \{1, \dots, d\}$ and therefore only a number linear in n of feasible index sets has to be checked. This is essentially the method proposed by [6]. The drawback of this approach is that the time complexity is quasilinear in n because of the sorting, and the space complexity is linear in n because the permutation has to be remembered to be undone afterwards.

4 Finding the Zero of the Auxiliary Function

Lemma 4 gives a compact expression that characterizes projections onto D . We first note that the representation only depends on one number:

Remark 5. Let $x \in \mathbb{R}_{\geq 0}^n \setminus D$ such that $\text{proj}_D(x)$ is unique. Then there is exactly one $\alpha \in \mathbb{R}$ such that $\text{proj}_D(x) = \frac{\lambda_2 \cdot \max(x - \alpha \cdot e, 0)}{\|\max(x - \alpha \cdot e, 0)\|_2}$.

Proof. The projection becomes $\text{proj}_D(x) = \max(\frac{1}{\beta}(x - \alpha \cdot e), 0)$ with unique numbers $\alpha \in \mathbb{R}$ and $\beta \in \mathbb{R}_{>0}$ due to Lemma 4. With $\beta > 0$ we have that $\lambda_2 = \|\max(\frac{1}{\beta}(x - \alpha \cdot e), 0)\|_2 = \frac{1}{\beta} \|\max(x - \alpha \cdot e, 0)\|_2$, and the claim follows. \square

It can hence be concluded that the sparseness projection can be considered a soft variant of thresholding [9]:

Definition 6. The function $\mathcal{S}_\alpha: \mathbb{R} \rightarrow \mathbb{R}$, $x \mapsto \max(x - \alpha, 0)$, is called *soft-shrinkage function*, where $\alpha \in \mathbb{R}$. It is continuous on \mathbb{R} and differentiable exactly on $\mathbb{R} \setminus \{\alpha\}$.

With Remark 5 we know that we only have to find one scalar to compute projections onto D . Analogous to the projection onto a simplex [7], we can thus define an auxiliary function which vanishes exactly at the number that yields the projection:

Definition 7. Let $x \in \mathbb{R}_{\geq 0}^n \setminus D$ such that $\text{proj}_D(x)$ is unique and $\sigma(x) < \sigma^*$. Let the maximum entry of x be denoted by $x_{\max} := \max_{i \in \{1, \dots, n\}} x_i$. Then the function

$$\Psi: [0, x_{\max}) \rightarrow \mathbb{R}, \quad \alpha \mapsto \frac{\|\max(x - \alpha \cdot e, 0)\|_1}{\|\max(x - \alpha \cdot e, 0)\|_2} - \frac{\lambda_1}{\lambda_2}$$

is called *auxiliary function* to the projection onto D .

Note that the case of $\sigma(x) \geq \sigma^*$ is trivial, because in this sparseness-decreasing setup we have that all coordinates of the projection must be positive. Hence $I = \{1, \dots, n\}$ in the proof of Lemma 4, and the shifting scalar α can be computed from a closed-form expression.

We further fix some notation for convenience:

Definition 8. Let $x \in \mathbb{R}_{\geq 0}^n$ be a vector. Then we write $\mathcal{X} := \{x_i \mid i \in \{1, \dots, n\}\}$ for the set of entries of x . Further, let $x_{\min} := \min \mathcal{X}$ be short for the smallest entry of x , and $x_{\max} := \max \mathcal{X}$ and $x_{2\text{nd-max}} := \max \mathcal{X} \setminus \{x_{\max}\}$ denote the two largest entries of x .

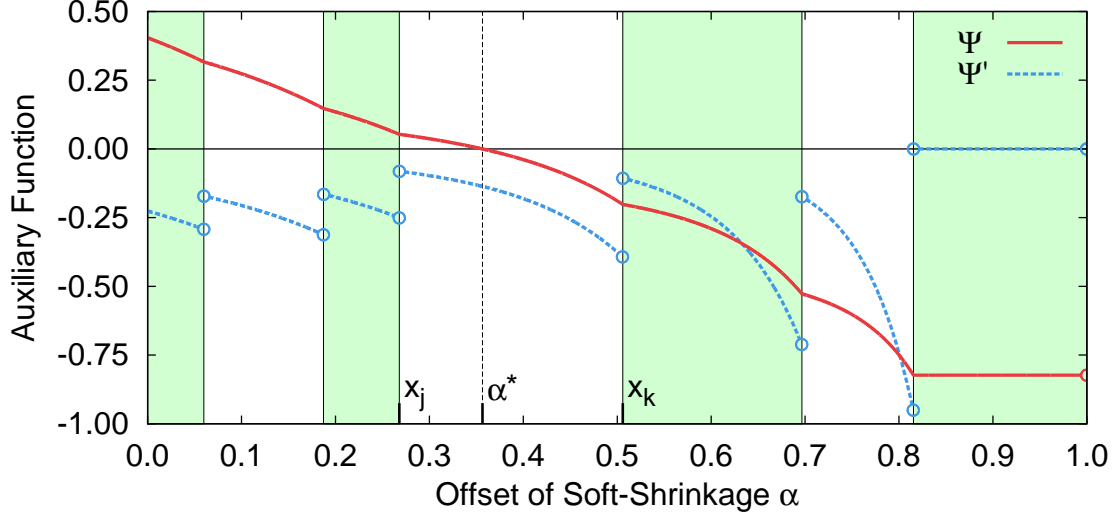


Figure 1: Plot of the auxiliary function Ψ and its derivative for a random vector x (see Lemma 9 for an analysis). The derivative Ψ' was scaled using a positive number for improved visibility. The steps in Ψ' are exactly the places where α coincides with an entry of x . With Remark 11, it is sufficient to find an α such that $\Psi(x_j) \geq 0$ and $\Psi(x_k) < 0$ for the neighboring entries x_j and x_k in x , because then the exact solution α^* can be computed with a closed-form expression.

Let $q: \mathbb{R} \rightarrow \mathbb{R}^n$, $\alpha \mapsto \max(x - \alpha \cdot e, 0)$, denote the curve that evolves from entry-wise application of the soft-shrinkage function. Let the Manhattan norm and Euclidean norm of points from q be given by $\ell_1: \mathbb{R} \rightarrow \mathbb{R}$, $\alpha \mapsto \|q(\alpha)\|_1$, and $\ell_2: \mathbb{R} \rightarrow \mathbb{R}$, $\alpha \mapsto \|q(\alpha)\|_2$, respectively. We thus have that $\Psi = \frac{\ell_1}{\ell_2} - \frac{\lambda_1}{\lambda_2}$, and we find that $q(\alpha) \neq 0$ if and only if $\alpha < x_{\max}$.

In order to efficiently find the zeros of Ψ we first investigate its analytical properties. See Figure 1 for an example of Ψ that provides orientation for the next result.

Lemma 9. Let $x \in \mathbb{R}_{\geq 0}^n \setminus D$ be given such that the auxiliary function Ψ is well-defined. Then:

- (a) Ψ is continuous on $[0, x_{\max})$.
- (b) Ψ is differentiable on $[0, x_{\max}) \setminus \mathcal{X}$.
- (c) Ψ is strictly decreasing on $[0, x_{2\text{nd-max}})$ and constant on $[x_{2\text{nd-max}}, x_{\max})$.
- (d) There is exactly one $\alpha^* \in (0, x_{2\text{nd-max}})$ with $\Psi(\alpha^*) = 0$.
- (e) $\text{proj}_D(x) = \frac{\lambda_2 \cdot \max(x - \alpha^* \cdot e, 0)}{\|\max(x - \alpha^* \cdot e, 0)\|_2}$ where α^* is the zero of Ψ .

Proof. (a) q is continuous because the soft-shrinkage function is continuous. Hence so are ℓ_1 and ℓ_2 , and hence Ψ as compositions of continuous functions.

(b) The soft-shrinkage function causes Ψ to be differentiable exactly on $[0, x_{\max}) \setminus \mathcal{X}$. Now let $x_j < x_k$ be two successive elements from \mathcal{X} , such that there is no element from \mathcal{X} between them. In the case that $x_k = x_{\min}$ it can be assumed that $x_j = 0$. Then the index set

$I := \{i \in \{1, \dots, n\} \mid x_i > \alpha\}$ of non-vanishing coordinates in q is constant for $\alpha \in (x_j, x_k)$, and the derivative of Ψ can be computed using a closed-form expression. For this, let $d := |I|$ denote the number of nonzero entries in q . With $\ell_1(\alpha) = \sum_{i \in I} (x_i - \alpha) = \sum_{i \in I} x_i - d\alpha$ we obtain $\ell'_1(\alpha) = -d$. Analogously, it is $\frac{\partial}{\partial \alpha} \ell_2(\alpha)^2 = \frac{\partial}{\partial \alpha} \sum_{i \in I} (x_i - \alpha)^2 = -2 \sum_{i \in I} (x_i - \alpha) = -2\ell_1(\alpha)$, and hence $\ell'_2(\alpha) = \frac{\partial}{\partial \alpha} \sqrt{\ell_2(\alpha)^2} = \frac{1}{2} \ell_2(\alpha)^{-1} \frac{\partial}{\partial \alpha} \ell_2(\alpha)^2 = -\frac{\ell_1(\alpha)}{\ell_2(\alpha)}$. Therefore, the quotient rule yields

$$\Psi'(\alpha) = \frac{-d\ell_2(\alpha) + \ell_1(\alpha) \frac{\ell'_1(\alpha)}{\ell_2(\alpha)}}{\ell_2(\alpha)^2} = \frac{1}{\ell_2(\alpha)} \left(\frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} - d \right).$$

It can further be shown that higher derivatives are of similar form. We have that $\frac{\partial}{\partial \alpha} \ell_1(\alpha)^2 = 2\ell_1(\alpha)\ell'_1(\alpha) = -2d\ell_1(\alpha)$, and thus

$$\frac{\partial}{\partial \alpha} \frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} = \frac{-2d\ell_1(\alpha)\ell_2(\alpha)^2 + 2\ell_1(\alpha)^3}{\ell_2(\alpha)^4} = 2 \frac{\ell_1(\alpha)}{\ell_2(\alpha)^2} \left(\frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} - d \right).$$

We also obtain $\frac{\partial}{\partial \alpha} \frac{1}{\ell_2(\alpha)} = \frac{-\ell'_2(\alpha)}{\ell_2(\alpha)^2} = \frac{\ell_1(\alpha)}{\ell_2(\alpha)^3}$, and eventually

$$\Psi''(\alpha) = \frac{\ell_1(\alpha)}{\ell_2(\alpha)^3} \left(\frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} - d \right) + \frac{2}{\ell_2(\alpha)} \frac{\ell_1(\alpha)}{\ell_2(\alpha)^2} \left(\frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} - d \right) = 3 \frac{\ell_1(\alpha)}{\ell_2(\alpha)^3} \left(\frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} - d \right),$$

or in other words $\frac{\Psi''(\alpha)}{\Psi'(\alpha)} = 3 \frac{\ell_1(\alpha)}{\ell_2(\alpha)^2}$.

(c) First let $\alpha \in (x_{2\text{nd-max}}, x_{\text{max}})$. With the notation of (b) we then have that $d = 1$, such that q has exactly one non-vanishing coordinate. Hence, $\ell_1(\alpha) = \ell_2(\alpha)$ and $\Psi' \equiv 0$ on $(x_{2\text{nd-max}}, x_{\text{max}})$, thus Ψ is constant on $(x_{2\text{nd-max}}, x_{\text{max}})$ as a consequence of the mean value theorem from real analysis. Because Ψ is continuous, it is constant even on $[x_{2\text{nd-max}}, x_{\text{max}})$.

Next let $\alpha \in [0, x_{2\text{nd-max}}) \setminus \mathcal{X}$, then $d \geq 2$ and $\ell_1(\alpha) \leq \sqrt{d}\ell_2(\alpha)$ with Remark 1. The inequality is in fact strict, because $q(\alpha)$ has at least two distinct nonzero entries. This implies that $\Psi' < 0$ on (x_j, x_k) where $x_j < x_k$ are neighbors of α as in (b). The mean value theorem then guarantees that Ψ is strictly decreasing between neighboring elements from \mathcal{X} . This property holds then for the entire interval $[0, x_{2\text{nd-max}})$ due to the continuity of Ψ .

(d) We have by requirement that $\sigma(x) < \sigma^*$, and therefore $\frac{\|x\|_1}{\|x\|_2} > \frac{\lambda_1}{\lambda_2}$, and so $\Psi(0) > 0$. For $\alpha \in (x_{2\text{nd-max}}, x_{\text{max}})$ we obtain $\ell_1(\alpha) = \ell_2(\alpha)$ as in (c). It is then $\Psi(\alpha) < 0$ using $\lambda_2 < \lambda_1$. The existence of $\alpha^* \in [0, x_{2\text{nd-max}})$ with $\Psi(\alpha^*) = 0$ follows from the intermediate value theorem and (c). Uniqueness of α^* is guaranteed because Ψ is strictly monotone.

(e) With Remark 5 there is exactly one $\tilde{\alpha} \in \mathbb{R}$ such that $\text{proj}_D(x) = \frac{\lambda_2 \cdot \max(x - \tilde{\alpha} \cdot e, 0)}{\|\max(x - \tilde{\alpha} \cdot e, 0)\|_2}$. We see that $\Psi(\tilde{\alpha}) = 0$, and the uniqueness of the zero of Ψ implies that $\alpha^* = \tilde{\alpha}$. \square

The unique zero of the auxiliary function can be found numerically using standard root-finding algorithms, such as Bisection, Newton's method or Halley's method [10]. We can improve on this by noting that whenever a number is found in a certain interval, then the exact value of the zero of Ψ can already be computed.

Remark 10. Let $x \in \mathbb{R}_{\geq 0}^n \setminus D$ such that the auxiliary function Ψ is well-defined. Then there are two unique numbers $x_j < x_k$, where either $x_j = 0$ and $x_k = x_{\text{min}}$ or $x_j, x_k \in \mathcal{X}$ such that there is no other element from \mathcal{X} in between, such that $\Psi(x_j) \geq 0$ and $\Psi(x_k) < 0$ and there is an $\alpha^* \in [x_j, x_k)$ with $\Psi(\alpha^*) = 0$.

Proof. Let $\alpha^* \in (0, x_{2\text{nd-max}})$ with $\Psi(\alpha^*) = 0$ be given with Lemma 9. When $\alpha^* < x_{\min}$ holds, existence follows immediately with Lemma 9 by setting $x_j := 0$ and $x_k := x_{\min}$. Otherwise, define $x_j := \max\{x_i \mid x_i \in \mathcal{X} \text{ and } x_i \leq \alpha^*\}$ and $x_k := \min\{x_i \mid x_i \in \mathcal{X} \text{ and } x_i > \alpha^*\}$, which both exist as the sets where the maximum and the minimum is taken are nonempty. Clearly these two numbers fulfill the condition from the claim by Lemma 9. The bracketing by x_j and x_k is unique because α^* is in both cases unique with Lemma 9. \square

We further note that it is easy to check whether the correct interval has already been found, and give a closed-form expression for the zero of Ψ in this case:

Remark 11. Let $x \in \mathbb{R}_{\geq 0}^n \setminus D$ such that the auxiliary function Ψ is well-defined and let $\alpha \in [0, x_{\max})$. If $\alpha < x_{\min}$ define $x_j := 0$ and $x_k := x_{\min}$, otherwise let $x_j \leq \alpha < x_k$ with $x_j, x_k \in \mathcal{X}$ such that there is no element from \mathcal{X} between x_j and x_k . Let $I := \{i \in \{1, \dots, n\} \mid x_i > \alpha\} = \{i_1, \dots, i_d\}$ where $d := |I|$ and $\tilde{x} \in \mathbb{R}_{\geq 0}^d$ such that $\tilde{x}^T = (x_{i_1}, \dots, x_{i_d})$. Then the following holds:

- (a) $\ell_1(\xi) = \|\tilde{x}\|_1 - d\xi$ and $\ell_2^2(\xi) = \|\tilde{x}\|_2^2 - 2\xi\|\tilde{x}\|_1 + d\xi^2$ for $\xi \in \{x_j, \alpha, x_k\}$.
- (b) When $\lambda_2 \ell_1(x_j) \geq \lambda_1 \ell_2(x_j)$ and $\lambda_2 \ell_1(x_k) < \lambda_1 \ell_2(x_k)$ hold, then

$$\alpha^* := \frac{1}{d} \left(\|\tilde{x}\|_1 - \lambda_1 \sqrt{\frac{d\|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2}{d\lambda_2^2 - \lambda_1^2}} \right)$$

is the unique zero of Ψ .

Proof. (a) We have that $\ell_1(\alpha) = \sum_{i \in I} (x_i - \alpha) = \sum_{i \in I} x_i - d\alpha = \|\tilde{x}\|_1 - d\alpha$ and further $\ell_2(\alpha)^2 = \sum_{i \in I} (x_i - \alpha)^2 = \sum_{i \in I} (x_i^2 - 2\alpha x_i + \alpha^2) = \|\tilde{x}\|_2^2 - 2\alpha\|\tilde{x}\|_1 + d\alpha^2$.

Now let $K := \{i \in \{1, \dots, n\} \mid x_i > x_k\}$ and $\tilde{K} := \{i \in \{1, \dots, n\} \mid x_i = x_k\}$. Then $K = I \setminus \tilde{K}$, and thus $\ell_1(x_k) = \sum_{i \in K} (x_i - x_k) = \sum_{i \in I} (x_i - x_k) - \sum_{i \in \tilde{K}} (x_i - x_k) = \sum_{i \in I} (x_i - x_k) = \|\tilde{x}\|_1 - dx_k$. Likewise follows $\ell_2(x_k)^2 = \|\tilde{x}\|_2^2 - 2x_k\|\tilde{x}\|_1 + dx_k^2$.

Finally, let $J := \{i \in \{1, \dots, n\} \mid x_i > x_j\}$ and $\tilde{J} := \{i \in \{1, \dots, n\} \mid x_i = x_j\}$. Then $I = J \setminus \tilde{J}$, and we obtain $\ell_1(x_j) = \sum_{i \in J} (x_i - x_j) = \sum_{i \in I} (x_i - x_j) + \sum_{i \in \tilde{J}} (x_i - x_j) = \sum_{i \in I} (x_i - x_j) = \|\tilde{x}\|_1 - dx_j$. The claim for $\ell_2(x_j)^2$ follows analogously.

(b) The condition from the claim is equivalent to $\Psi(x_j) \geq 0$ and $\Psi(x_k) < 0$. Hence with Remark 10 there is an $\alpha^* \in [x_j, x_k)$ with $\Psi(\alpha^*) = 0$. Let $p := \text{proj}_D(x)$ be the projection of x onto D and define $J := \{i \in \{1, \dots, n\} \mid p_i > 0\}$. Lemma 9(e) implies $J = \{i \in \{1, \dots, n\} \mid x_i > \alpha^*\}$. Furthermore, it is $J = \{i \in \{1, \dots, n\} \mid x_i > x_j\} = \{i \in \{1, \dots, n\} \mid x_i > \alpha\} = I$. Thus we already had the correct set of non-vanishing coordinates of the projection in the first place, and the expression for α^* follows from the proof of Lemma 4. \square

5 A Linear Time and Constant Space Projection Algorithm

By exploiting the analytical properties of Ψ , simple methods are sufficient to locate the interval in which its zero resides. Because the interval has a positive length, simple Bisection is guaranteed to find it in a constant number of steps [7]. Empirically, we found that solvers that use the

Algorithm 1: Linear time and constant space evaluation of the auxiliary function Ψ .

Input: $x \in \mathbb{R}_{\geq 0}^n$, $\lambda_1, \lambda_2, \alpha \in \mathbb{R}$ with $0 < \lambda_2 < \lambda_1 < \sqrt{n}\lambda_2$ and $0 \leq \alpha < \max_{i \in \{1, \dots, n\}} x_i$.
Output: $\Psi(\alpha), \Psi'(\alpha), \Psi''(\alpha), \tilde{\Psi}(\alpha), \tilde{\Psi}'(\alpha) \in \mathbb{R}$, $\text{finished} \in \mathbb{B}$, $\ell_1, \ell_2^2 \in \mathbb{R}$, $d \in \mathbb{N}$.

```
// Initialize.
1  $\ell_1 := 0$ ;  $\ell_2^2 := 0$ ;  $d := 0$ ;  $x_j := 0$ ;  $\Delta x_j := -\alpha$ ;  $x_k := \infty$ ;  $\Delta x_k := \infty$ ;
   // Scan through  $x$ .
2 for  $i := 1$  to  $n$  do
3    $t := x_i - \alpha$ ;
4   if  $t > 0$  then
5      $\ell_1 := \ell_1 + x_i$ ;  $\ell_2^2 := \ell_2^2 + x_i^2$ ;  $d := d + 1$ ;
6     if  $t < \Delta x_k$  then  $x_k := x_i$ ;  $\Delta x_k := t$ ; ;
7   else
8     if  $t > \Delta x_j$  then  $x_j := x_i$ ;  $\Delta x_j := t$ ; ;
9   end
10 end

   // Compute  $\Psi(\alpha)$ ,  $\Psi'(\alpha)$  and  $\Psi''(\alpha)$ .
11  $\ell_1(\alpha) := \ell_1 - d\alpha$ ;  $\ell_2(\alpha)^2 := \ell_2^2 - 2\alpha\ell_1 + d\alpha^2$ ;
12  $\Psi(\alpha) := \frac{\ell_1(\alpha)}{\sqrt{\ell_2(\alpha)^2}} - \frac{\lambda_1}{\lambda_2}$ ;  $\Psi'(\alpha) := \frac{1}{\sqrt{\ell_2(\alpha)^2}} \left( \frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} - d \right)$ ;  $\Psi''(\alpha) := \frac{3\Psi'(\alpha)\ell_1(\alpha)}{\ell_2(\alpha)^2}$ ;

   // Compute  $\tilde{\Psi}(\alpha)$  and  $\tilde{\Psi}'(\alpha)$ .
13  $\tilde{\Psi}(\alpha) := \frac{\ell_1(\alpha)^2}{\ell_2(\alpha)^2} - \frac{\lambda_1^2}{\lambda_2^2}$ ;  $\tilde{\Psi}'(\alpha) := 2 \frac{\ell_1(\alpha)}{\sqrt{\ell_2(\alpha)^2}} \Psi'(\alpha)$ ;

   // Compute  $\Psi(x_j)$  and  $\Psi(x_k)$ , check for sign change and return.
14  $\text{finished} := \lambda_2(\ell_1 - dx_j) \geq \lambda_1 \sqrt{\ell_2^2 - 2x_j\ell_1 + dx_j^2}$  and  $\lambda_2(\ell_1 - dx_k) < \lambda_1 \sqrt{\ell_2^2 - 2x_k\ell_1 + dx_k^2}$ ;
15 return  $(\Psi(\alpha), \Psi'(\alpha), \Psi''(\alpha), \tilde{\Psi}(\alpha), \tilde{\Psi}'(\alpha), \text{finished}, \ell_1, \ell_2^2, d)$ ;
```

derivative of Ψ converge faster, despite of the step discontinuities of Ψ' . We have implemented Newton's method, Halley's method, and Newton's method applied to the slightly transformed auxiliary function $\tilde{\Psi} := \frac{\ell_1^2}{\ell_2^2} - \frac{\lambda_1^2}{\lambda_2^2}$. These methods were additionally safeguarded with Bisection to guarantee new positions are located within well-defined bounds [11]. This does impair the theoretical property that only a constant number of steps be required to find a solution, but in practice a significantly smaller number of steps needs to be made compared to plain Bisection. This is demonstrated through experimental results in Section 7.

We are now in a position to formulate the main result of this paper, by proposing an efficient algorithm for computing sparseness-enforcing projections:

Theorem 12. Algorithm 2 computes projections onto D , where unique, in a number of operations linear in the problem dimensionality n and with only constant additional space.

The proof is omitted as it is essentially a composition of the results from Section 4.

Algorithm 2: Linear time and constant space projection onto D . The auxiliary function Ψ is evaluated by calls to "auxiliary", which are carried out by Algorithm 1.

Input: $x \in \mathbb{R}_{\geq 0}^n, \lambda_1, \lambda_2 \in \mathbb{R}$ with $0 < \lambda_2 < \lambda_1 < \sqrt{n}\lambda_2$,
 solver $\in \{\text{Bisection}, \text{Newton}, \text{NewtonSqr}, \text{Halley}\}$.

Output: $\text{proj}_D(x) \in D$ where $D := S_{\geq 0}^{(\lambda_1, \lambda_2)} \subseteq \mathbb{R}_{\geq 0}^n$.

```

// Check whether sparseness should be increased or decreased.
1  $(\Psi(\alpha), \Psi'(\alpha), \Psi''(\alpha), \tilde{\Psi}(\alpha), \tilde{\Psi}'(\alpha), \text{finished}, \ell_1, \ell_2^2, d) := \text{auxiliary}(x, \lambda_1, \lambda_2, 0);$ 
2 if  $\Psi(\alpha) \leq 0$  then go to Line 18; // Decrease sparseness, skip root-finding.

// Need to increase sparseness, initialize safeguarded root-finding.
3  $\text{lo} := 0; \text{up} := \max\{x_i \mid i \in \{1, \dots, n\}, x_i \neq \max_{j \in \{1, \dots, n\}} x_j\}; \alpha := \text{lo} + \frac{1}{2}(\text{up} - \text{lo});$ 
4  $(\Psi(\alpha), \Psi'(\alpha), \Psi''(\alpha), \tilde{\Psi}(\alpha), \tilde{\Psi}'(\alpha), \text{finished}, \ell_1, \ell_2^2, d) := \text{auxiliary}(x, \lambda_1, \lambda_2, \alpha);$ 

// Perform root-finding until correct interval has been found.
5 while not finished do
    // Update Bisection interval.
    6 if  $\Psi(\alpha) > 0$  then  $\text{lo} := \alpha$  else  $\text{up} := \alpha$ ;
    // One iteration of root-finding.
    7 if solver = Bisection then  $\alpha := \text{lo} + \frac{1}{2}(\text{up} - \text{lo});$ 
    8 else // Use solvers based on derivatives.
        9 if solver = Newton then  $\alpha := \alpha - \frac{\Psi(\alpha)}{\Psi'(\alpha)};$ 
        10 else if solver = NewtonSqr then  $\alpha := \alpha - \frac{\tilde{\Psi}(\alpha)}{\tilde{\Psi}'(\alpha)};$ 
        11 else if solver = Halley then
            12  $h := 1 - \frac{\Psi(\alpha)\Psi''(\alpha)}{2\Psi'(\alpha)^2}; h := \max(0.5, \min(1.5, h)); \alpha := \alpha - \frac{\Psi(\alpha)}{h\Psi'(\alpha)};$ 
        13 end
        // If  $\alpha$  fell out of bounds, perform normal Bisection.
        14 if  $\alpha < \text{lo}$  or  $\alpha > \text{up}$  then  $\alpha := \text{lo} + \frac{1}{2}(\text{up} - \text{lo});$ 
    15 end
    // Re-evaluate auxiliary function at new position.
    16  $(\Psi(\alpha), \Psi'(\alpha), \Psi''(\alpha), \tilde{\Psi}(\alpha), \tilde{\Psi}'(\alpha), \text{finished}, \ell_1, \ell_2^2, d) := \text{auxiliary}(x, \lambda_1, \lambda_2, \alpha);$ 
17 end

// Correct interval has been found, compute exact value for  $\alpha$ .
18  $\alpha := \frac{1}{d} \left( \ell_1 - \lambda_1 \sqrt{\frac{d\ell_2^2 - \ell_1^2}{d\lambda_2^2 - \lambda_1^2}} \right);$ 

// Compute result of the projection in-place.
19  $\rho := 0;$ 
20 for  $i := 1$  to  $n$  do
    21  $t := x_i - \alpha;$ 
    22 if  $t > 0$  then  $x_i := t; \rho := \rho + t^2$  else  $x_i := 0;$ 
23 end
24 for  $i := 1$  to  $n$  do  $x_i := \frac{\lambda_2}{\sqrt{\rho}} x_i;$ 
25 return  $x;$ 
```

6 Gradient of the Projection

We conclude our analysis of the sparseness-enforcing projection by considering its gradient:

Lemma 13. The projection onto D can be cast as function $\mathbb{R}_{\geq 0}^n \rightarrow D$ in all points with unique projections, that is almost everywhere. Further, this function is differentiable almost everywhere.

More precisely, let $x \in \mathbb{R}_{\geq 0}^n \setminus D$ such that $p := \text{proj}_D(x)$ is unique. With Remark 5, let $\alpha \in \mathbb{R}$ such that $p = \frac{\lambda_2 \cdot \max(x - \alpha \cdot e, 0)}{\|\max(x - \alpha \cdot e, 0)\|_2}$. If $x_i \neq \alpha$ for all $i \in \{1, \dots, n\}$, then proj_D is differentiable in x . It is further possible to give a closed-form expression for the gradient as follows. Let the index set of nonzero entries in the projection be denoted by $I := \{i \in \{1, \dots, n\} \mid p_i > 0\} = \{i_1, \dots, i_d\}$ where $d := |I|$. Let $e_k \in \mathbb{R}^n$ denote the k th canonical basis vector for $k \in \{1, \dots, n\}$ and let $V := (e_{i_1}, \dots, e_{i_d})^T \in \{0, 1\}^{d \times n}$ be the slicing matrix with respect to I , that is with $\tilde{x} := Vx \in \mathbb{R}^d$ we have for example $\tilde{x}^T = (x_{i_1}, \dots, x_{i_d})$. Write $a := d\|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2 \in \mathbb{R}_{\geq 0}$ and $b := d\lambda_2^2 - \lambda_1^2 \in \mathbb{R}_{\geq 0}$ for short. With Lemma 4 we find that $\alpha = \frac{1}{d}(\|\tilde{x}\|_1 - \lambda_1 \sqrt{\frac{a}{b}})$. Denote by $\tilde{e} := Ve \in \{1\}^d$ the vector where all d entries are equal to unity, and let $\tilde{q} := \max(\tilde{x} - \alpha \cdot \tilde{e}, 0) = \tilde{x} - \alpha \cdot \tilde{e} \in \mathbb{R}_{\geq 0}^d$ which implies that $p = \frac{\lambda_2}{\|\tilde{q}\|_2} V^T \tilde{q}$ holds.

Let $\tilde{p} := \frac{\lambda_2}{\|\tilde{q}\|_2} \tilde{q}$ such that $p = V^T \tilde{p}$, then $\frac{\partial}{\partial x} \text{proj}_D(x) = V^T G V \in \mathbb{R}^{n \times n}$ where

$$G := \sqrt{\frac{b}{a}} E_d - \frac{1}{\sqrt{ab}} (\lambda_2^2 \tilde{e} \tilde{e}^T + d \tilde{p} \tilde{p}^T - \lambda_1 (\tilde{e} \tilde{p}^T + \tilde{p} \tilde{e}^T)).$$

Proof. The projection is unique almost everywhere as already shown by [5]. When $x_i \neq \alpha$ for all $i \in \{1, \dots, n\}$, proj_D is differentiable as composition of differentiable functions as then I is invariant to local changes in x . Write $\tilde{p} := \frac{\lambda_2}{\|\tilde{q}\|_2} \tilde{q}$ such that $p = V^T \tilde{p}$, then the chain rule yields

$$\frac{\partial p}{\partial x} = \frac{\partial V^T \tilde{p}}{\partial \tilde{p}} \cdot \frac{\partial}{\partial \tilde{q}} \left(\frac{\lambda_2}{\|\tilde{q}\|_2} \tilde{q} \right) \cdot \frac{\partial (\tilde{x} - \alpha \cdot \tilde{e})}{\partial \tilde{x}} \cdot \frac{\partial Vx}{\partial x} = V^T \cdot \underbrace{\lambda_2 \frac{\partial}{\partial \tilde{q}} \left(\frac{\tilde{q}}{\|\tilde{q}\|_2} \right) \cdot \left(E_d - \tilde{e} \frac{\partial \alpha}{\partial \tilde{x}} \right)}_{=: G \in \mathbb{R}^{d \times d}} \cdot V.$$

We thus only have to show that the matrix G defined here matches the matrix from the claim. It is easy to see that the mapping from a vector to its normalized version has a simple gradient, that is we have that $\frac{\partial}{\partial \tilde{q}} \frac{\tilde{q}}{\|\tilde{q}\|_2} = \frac{1}{\|\tilde{q}\|_2} \left(E_d - \frac{\tilde{q} \tilde{q}^T}{\|\tilde{q}\|_2^2} \right)$. Because \tilde{q} and \tilde{x} have only non-negative entries, the canonical dot product with \tilde{e} yields essentially their L_1 norms. We hence obtain $\|\tilde{q}\|_1 = \tilde{e}^T \tilde{x} - \alpha \tilde{e}^T \tilde{e} = \lambda_1 \sqrt{\frac{a}{b}}$. Likewise, the L_2 norm of \tilde{q} equals

$$\begin{aligned} \|\tilde{q}\|_2^2 &= \|\tilde{x}\|_2^2 - 2\alpha \|\tilde{x}\|_1 + d\alpha^2 = \|\tilde{x}\|_2^2 - \alpha \left(\|\tilde{x}\|_1 + \lambda_1 \sqrt{\frac{a}{b}} \right) = \|\tilde{x}\|_2^2 - \frac{1}{d} \left(\|\tilde{x}\|_1^2 - \lambda_1^2 \frac{a}{b} \right) \\ &= \frac{1}{d} (d\|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2) + \lambda_1^2 \frac{a}{b} = \frac{a}{d} \left(1 + \frac{\lambda_1^2}{b} \right) = \lambda_2^2 \frac{a}{b}. \end{aligned}$$

To compute the gradient of α , we first note that b does not depend on \tilde{x} but a does. It is $\frac{\partial}{\partial \tilde{x}} a = 2d\tilde{x}^T - 2\|\tilde{x}\|_1 \tilde{e}^T \in \mathbb{R}^{1 \times d}$, and hence $\frac{\partial}{\partial \tilde{x}} \sqrt{a} = \frac{1}{\sqrt{a}} (d\tilde{x}^T - \|\tilde{x}\|_1 \tilde{e}^T) \in \mathbb{R}^{1 \times d}$. With $\tilde{x} = \tilde{q} + \alpha \cdot \tilde{e}$ follows $d\tilde{x} - \|\tilde{x}\|_1 \tilde{e} = d\tilde{q} - \lambda_1 \sqrt{\frac{a}{b}} \tilde{e}$, and hence

$$\frac{\partial}{\partial \tilde{x}} \alpha = \frac{1}{d} \tilde{e}^T - \frac{\lambda_1}{d\sqrt{ab}} (d\tilde{x}^T - \|\tilde{x}\|_1 \tilde{e}^T) = \left(\frac{1}{d} + \frac{\lambda_1^2}{db} \right) \tilde{e}^T - \frac{\lambda_1}{\sqrt{ab}} \tilde{q}^T = \frac{\lambda_2^2}{b} \tilde{e}^T - \frac{\lambda_1}{\sqrt{ab}} \tilde{q}^T \in \mathbb{R}^{1 \times d}.$$

Algorithm 3: Product of the gradient of the projection onto D with an arbitrary vector.

Input: $y \in \mathbb{R}^n$ and the following results of Algorithm 2: $I = \{i_1, \dots, i_d\} \subseteq \{1, \dots, n\}$,
 $d := |I|$, $\tilde{p} \in \mathbb{R}_{\geq 0}^d$, $\lambda_1, \lambda_2 \in \mathbb{R}_{> 0}$, $a := d \|\tilde{x}\|_2^2 - \|\tilde{x}\|_1^2 \in \mathbb{R}_{\geq 0}$ and $b := d\lambda_2^2 - \lambda_1^2 \in \mathbb{R}_{\geq 0}$.

Output: $z := \left(\frac{\partial}{\partial x} \text{proj}_D(x) \right) \cdot y \in \mathbb{R}^n$.

```

// Scan and slice input vector.
1  $\tilde{y} \in \{0\}^d$ ;  $\text{sum}_{\tilde{y}} := 0$ ;  $\text{scp}_{p,y} := 0$ ;
2 for  $i := 1$  to  $d$  do  $\text{sum}_{\tilde{y}} := \text{sum}_{\tilde{y}} + z_{i_j}$ ;  $\text{scp}_{p,y} := \text{scp}_{p,y} + \tilde{p}_j \cdot z_{i_j}$ ;  $\tilde{z}_j := z_{i_j}$ ;

// Compute product with gradient in sliced space.
3  $\tilde{z} := \sqrt{\frac{b}{a}} \tilde{z}$ ;  $\tilde{z} := \tilde{z} + \frac{1}{\sqrt{ab}} (\lambda_1 \cdot \text{sum}_{\tilde{y}} - d \cdot \text{scp}_{p,y}) \tilde{p}$ ;  $\tilde{z} := \tilde{z} + \frac{1}{\sqrt{ab}} (\lambda_1 \cdot \text{scp}_{p,y} - \lambda_2^2 \cdot \text{sum}_{\tilde{y}}) \tilde{e}$ ;

// Unslicing to yield final result.
4  $y \in \{0\}^n$ ; for  $i := 1$  to  $d$  do  $y_{i_j} := \tilde{z}_j$ ;
5 return  $y$ ;
```

Therefore, by substitution into G and multiplying out we yield

$$\begin{aligned}
G &= \sqrt{\frac{b}{a}} \left(E_d - \frac{b}{\lambda_2^2 a} \tilde{q} \tilde{q}^T \right) \left(E_d - \frac{\lambda_2^2}{b} \tilde{e} \tilde{e}^T + \frac{\lambda_1}{\sqrt{ab}} \tilde{e} \tilde{q}^T \right) \\
&= \sqrt{\frac{b}{a}} \left(E_d - \frac{\lambda_2^2}{b} \tilde{e} \tilde{e}^T + \frac{\lambda_1}{\sqrt{ab}} \tilde{e} \tilde{q}^T - \frac{b}{\lambda_2^2 a} \tilde{q} \tilde{q}^T + \frac{\lambda_1}{\sqrt{ab}} \tilde{q} \tilde{e}^T - \frac{\lambda_1^2}{\lambda_2^2 a} \tilde{q} \tilde{q}^T \right),
\end{aligned}$$

where we have used $\tilde{q} \tilde{q}^T \tilde{e} \tilde{e}^T = \tilde{q} (\tilde{q}^T \tilde{e}) \tilde{e}^T = \lambda_1 \sqrt{\frac{a}{b}} \tilde{q} \tilde{e}^T$ and $\tilde{q} \tilde{q}^T \tilde{e} \tilde{q}^T = \lambda_1 \sqrt{\frac{a}{b}} \tilde{q} \tilde{q}^T$. The claim then follows with $\frac{b}{\lambda_2^2 a} + \frac{\lambda_1^2}{\lambda_2^2 a} = \frac{d}{a}$ and $\tilde{q} = \sqrt{\frac{a}{b}} \tilde{p}$. \square

The gradient given in Lemma 13 has a particular simple form, as it is essentially a scaled identity matrix with additive combination of scaled dyadic products of simple vectors. In the situation where not the entire gradient but merely its product with an arbitrary vector is required, simple vector operations are already enough to compute the product:

Theorem 14. Algorithm 3 computes the product of the gradient of the sparseness projection with an arbitrary vector in time and space linear in the problem dimensionality n .

This claim can directly be validated using the expression from the gradient given in Lemma 13.

7 Experiments

To assess the performance of the algorithm we proposed to compute sparseness-enforcing projections, several experiments have been carried out. As the projection onto D is unique almost everywhere, different approaches must compute the same result except for a null set. We have compared the results of the algorithm proposed by [2] with the results of our algorithm for problem dimensionalities $n \in \{2^2, \dots, 2^{26}\}$ and for target degrees of sparseness $\sigma^* \in \{0.025, 0.050, \dots, 0.950, 0.975\}$. For every combination of n and σ^* we have sampled

one thousand random vectors, carried out both algorithms, and found that both algorithms produce numerically equal results given the very same input vector. Moreover, we have numerically verified the gradient of the projection for the same range using the central difference quotient.

Finally, experiments have been conducted to evaluate the choice of the solver for Algorithm 2. We have set the problem dimensionality to $n := 1024$, and then sampled one thousand random vectors for target sparseness degrees of $\sigma^* \in \{0.200, 0.225, \dots, 0.950, 0.975\}$. We have used the very same random vectors as input for all solvers, and counted the number of times the auxiliary function had to be computed until the solution was found. The results of this experiment are depicted in Figure 2. While Bisection needs about the same number of evaluations over all sparseness degrees, the solvers based on the derivative of Ψ depend on σ^* in their number of evaluations. This is because their starting value is set to the midpoint of the initial bracket in Algorithm 2, and thus their distance to the root of Ψ naturally depends on σ^* . The solver that performs best is NewtonSqr, that is Newton’s method applied to $\tilde{\Psi}$. It is quite surprising that the methods based on derivatives perform so well, as Ψ' possesses several step discontinuities as illustrated in Figure 1.

In the next experiment, the target sparseness degree was set to $\sigma^* := 0.90$ and the problem dimensionality n was varied in $\{2^2, \dots, 2^{26}\}$. The results are shown in Figure 3. The number of evaluations Bisection needs in the experiment grows about linearly in $\log(n)$. Because the expected minimum difference of two distinct entries from a random vector gets smaller when the dimensionality of the random vector is increased, the expected number of function evaluations Bisection requires increases with problem dimensionality. In either case, the length of the interval that has to be found is always bounded from below by the machine precision such that the number of function evaluations with Bisection is bounded from above. The methods based on derivatives exhibit sublinear growth, where the solver NewtonSqr is again the best performing one. Note that the number of iterations it requires decreases when dimensionality is enhanced. This is because Hoyer’s sparseness measure σ is not invariant to problem dimensionality, and hence a sparseness of $\sigma^* = 0.90$ has a different notion for $n = 2^{26}$ than for $n = 2^8$.

8 Conclusion

In this paper, we have proposed an efficient algorithm for computing sparseness-enforcing projections with respect to Hoyer’s sparseness measure σ . Although the target set of the projection is here non-convex, methods from projections onto simplexes could be adapted in a straightforward way. We have rigorously proved the correctness of our proposed algorithm, and additionally we have yielded a simple procedure to compute its gradient. We have shown that our algorithm needs only little resources, and that it scales well with problem dimensionality, even for very high target sparseness degrees.

Acknowledgments

The authors would like to thank Michael Gabb for helpful discussions. This work was supported by Daimler AG, Germany.

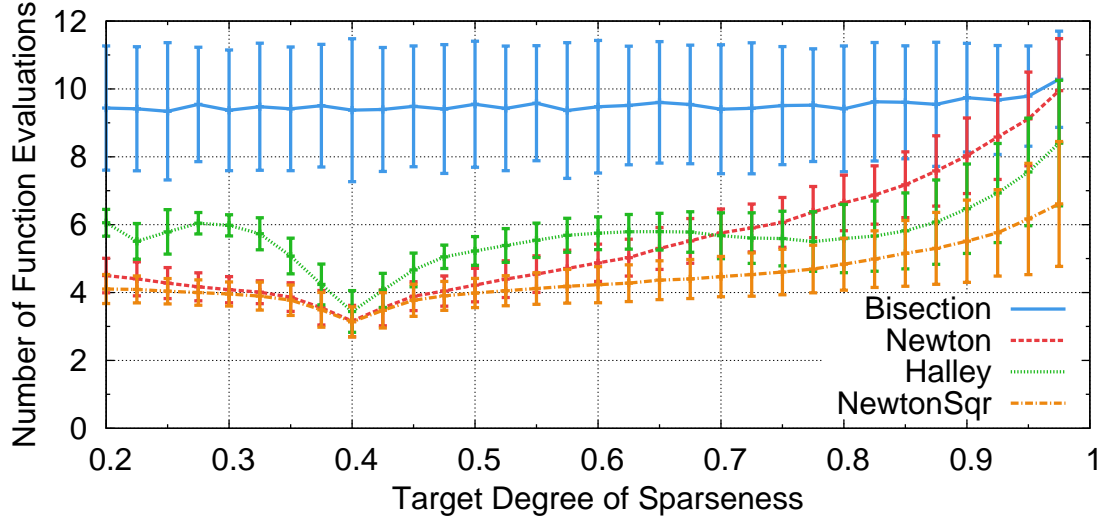


Figure 2: Auxiliary function evaluations needed to find the final interval with four different solvers. The problem dimensionality was set to $n := 1024$ and the target degree of sparseness σ^* was varied. While the performance of Bisection is constant over different values of σ^* , the solvers that use the derivative of the auxiliary function depend on the target sparseness and consistently outperform Bisection. Newton's method applied to $\tilde{\Psi}$ is the best-performing solver over all choices of σ^* .

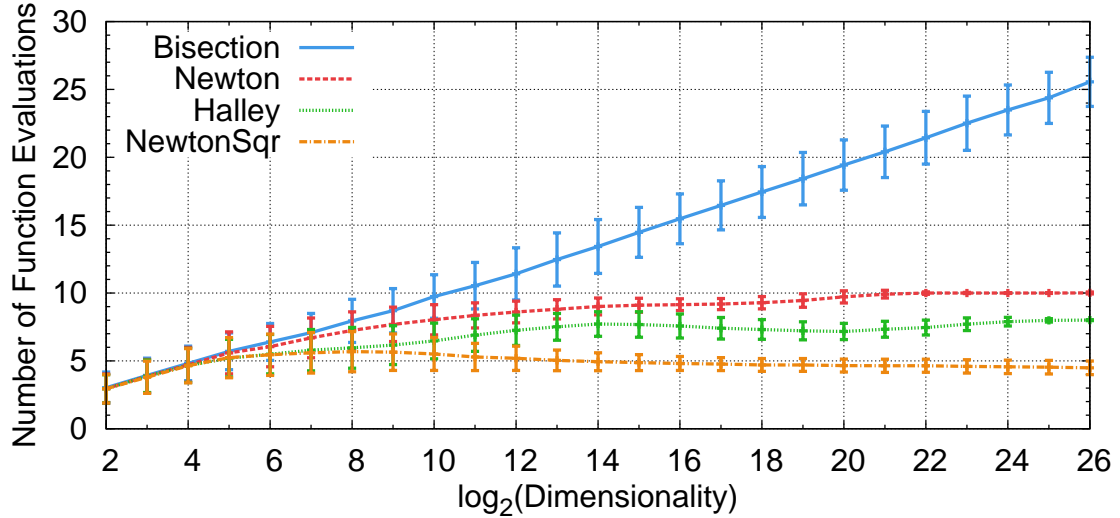


Figure 3: Same plot as in Figure 2, except for the target degree of sparseness was set to $\sigma^* := 0.90$ and the problem dimensionality was varied. The number of required function evaluations grows linearly with the logarithm of the problem dimensionality for Bisection, while the other solvers require a number sublinear in $\log(n)$. When $n = 2^{26} \approx 67 \cdot 10^6$ then Newton's method only requires 10 iterations in the mean, and Newton's method applied to $\tilde{\Psi}$ requires only 4 iterations.

References

- [1] N. Hurley and S. Rickard, “Comparing measures of sparsity,” *IEEE Transactions on Information Theory*, vol. 55, no. 10, pp. 4723–4741, 2009.
- [2] P. O. Hoyer, “Non-negative matrix factorization with sparseness constraints,” *Journal of Machine Learning Research*, vol. 5, pp. 1457–1469, 2004.
- [3] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Athena Scientific, 1999.
- [4] F. Deutsch, *Best Approximation in Inner Product Spaces*. Springer, 2001.
- [5] F. J. Theis, K. Stadlthanner, and T. Tanaka, “First results on uniqueness of sparse non-negative matrix factorization,” in *Proceedings of the European Signal Processing Conference*, 2005, pp. 1672–1675.
- [6] V. K. Potluru, S. M. Plis, J. L. Roux, B. A. Pearlmutter, V. D. Calhoun, and T. P. Hayes, “Block coordinate descent for sparse NMF,” Tech. Rep. arXiv:1301.3527v1, 2013.
- [7] J. Liu and J. Ye, “Efficient euclidean projections in linear time,” in *Proceedings of the International Conference on Machine Learning*, 2009, pp. 657–664.
- [8] A. J. Laub, *Matrix Analysis for Scientists and Engineers*. Society for Industrial and Applied Mathematics, 2004.
- [9] A. Hyvärinen, P. Hoyer, and E. Oja, “Sparse code shrinkage: Denoising by nonlinear maximum likelihood estimation,” in *Advances in Neural Information Processing Systems 11*, 1999, pp. 473–478.
- [10] J. F. Traub, *Iterative Methods for the Solution of Equations*. Prentice-Hall, 1964.
- [11] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes: The Art of Scientific Computing*, 3rd ed. Cambridge University Press, 2007.